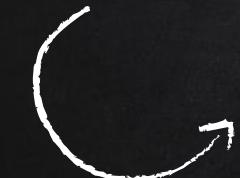


BASIC STATISTIC

--CONCEPTS/KNOWLEDGE/PRACTICE



主讲人:Lucy Xu

2021

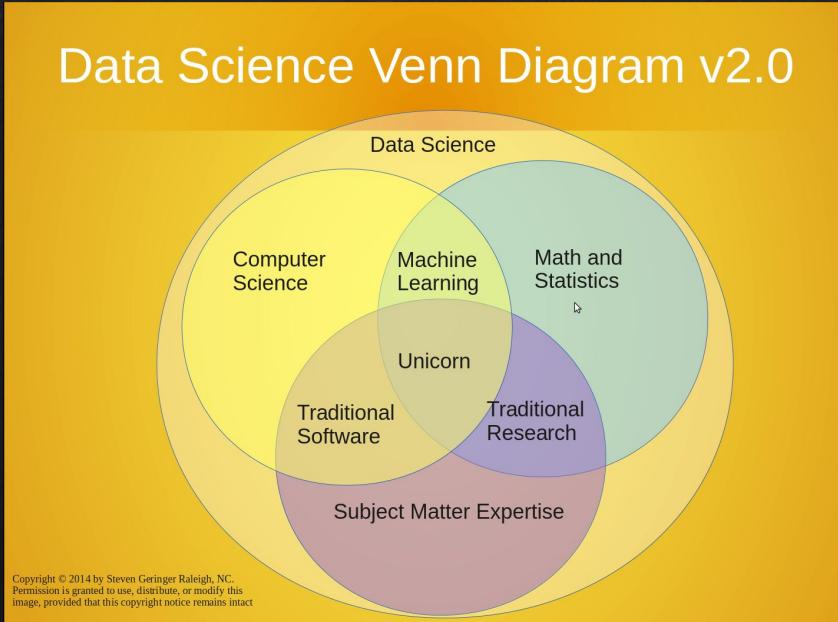


主讲人介绍

- 10+年数据科学, 预测建模等领域工作经验, 在消费者贷款领域拥有丰富的经验及知识, 包括房贷, 车贷, 信用卡和个人贷款。
- 现任全美前五保险公司, 机器学习和人工智能部门 – DATA SCIENCE DIRECTOR, 专注于IoT(INTERNET OF THINGS, 物联网) FOR CONNECTED CAR/HOME INSURANCE PRODUCT 和 NLP(NATURAL LANGUAGE PROCESSING)
- 曾任职CAPITAL ONE数据科学经理和DCF GLOBAL – SENIOR DATA SCIENCE MANAGER. 带领团队使用大数据和机器学习, 为整个信贷业务领域及产业链提供服务及解决方案。专注于市场营销, 信用风险, 定价, 估值, 运营, 帐户管理和欺诈检测。
- 本科毕业于中国科学技术大学, 在美国分别获得生物物理的博士和统计硕士学位



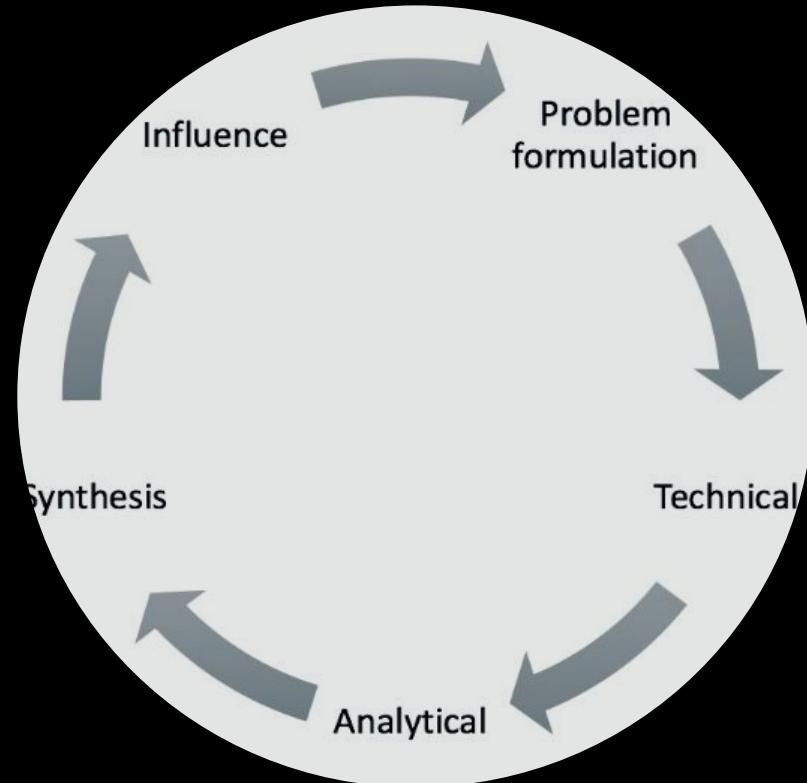
STATISTICS AND MATH ARE HEAVILY WEIGHTED FOR DATA SCIENTIST SKILL SETS



Battle of Data Science definition

“Data Scientist is a person who is better at statistics than any programmer and better at programming than any statistician.”

FIVE CORE SKILLS DATA SCIENTIST MANAGER IS LOOKING FOR





STATISTICS AND PROBABILITY ARE BUILDING BLOCKS OF ALL DATA HEAVY PRACTICE

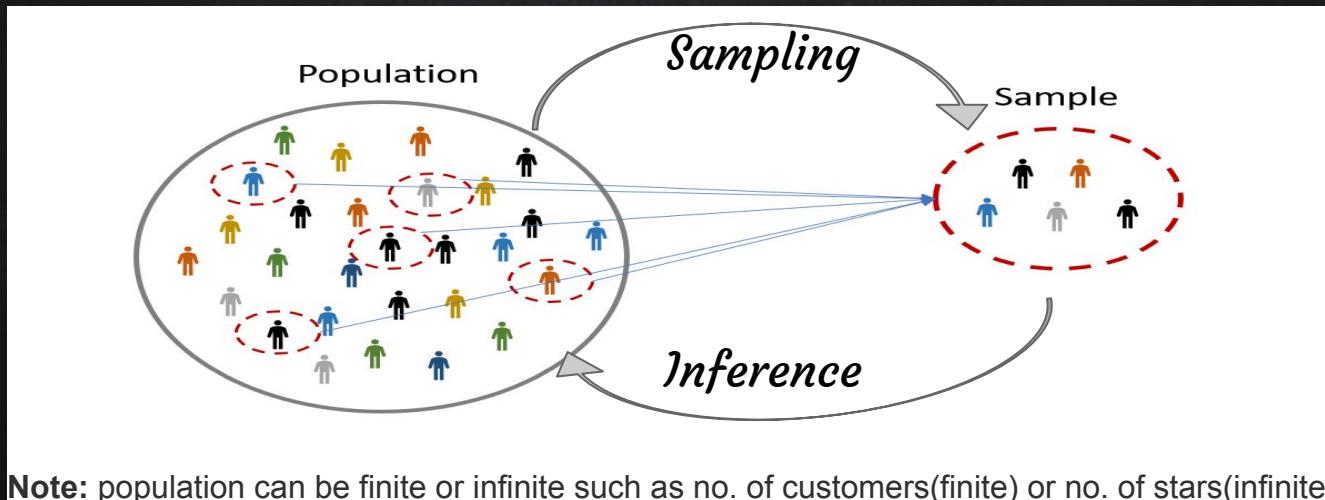


“Statistics is a mathematical science pertaining to the collection, analysis, interpretation or explanation, and presentation of data”

START WITH SOME BASIC TERMINOLOGIES – 1

Population: “The entire set of individuals or objects of interest”

Sample: “A portion or a part of the population of interest.”



Note: population can be finite or infinite such as no. of customers(finite) or no. of stars(infinite)



START WITH SOME BASIC TERMINOLOGIES – 2

Variable: “A characteristic of each individual element of a population/sample”.

Data: “A value of the variable associated with one element of a population or sample.it can be numbers, words or symbols.”

Experiment: “A planned activity whose result yields a set of data.”

Parameter: “A numerical value which summarizes the entire population or the descriptive measure of the population. For example, population mean, population variance, population standard deviation, etc.”

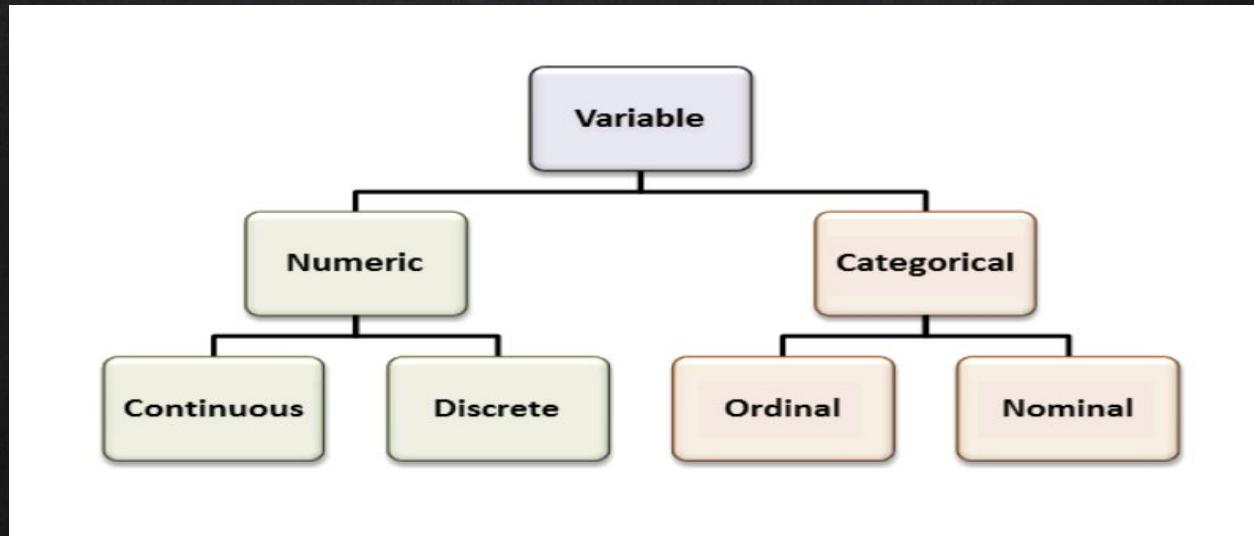
Statistic: “A numerical value which summarizes the sample data or descriptive measure of the sample. For example, sample mean, sample variance, sample standard deviation, etc. (Descriptive statistics and Inference statistics)”



TABLE OF CONTENTS FOR TODAY'S DISCUSSION

1. Basics Statistical concepts/methods to understand data (EDA statistics)
2. Sampling Techniques
3. Probability distribution and Inferential statistical analysis

VARIABLE/DATA TYPES



Independent(Features) variable and
dependent(Target/Outcome) variable



PRACTICE

Please label the data type, quantity/quality and continuous, for the list of variables (2-3 mins practice)

Sex,

State,

Grades(A/B/C/D),

number of children in a household,

monthly income,

Product rating (not recommend, neutral, may recommend, strongly recommend)

Company expenses

Customer ID



PRACTICE – ANSWER

Please label the data type, quantity/quality and continuous, for the list of variables (2–3 mins practice)

Sex – qualitative, Nominal

State – qualitative, Nominal

Grades(A/B/C/D) – qualitative, ordinal

number of children in a household – quantitative, discrete

monthly income – quantitative, continuous

Product rating (not recommend, neutral, may recommend, strongly recommend) – qualitative, ordinal

Company expenses – quantitative, continuous

Customer ID – qualitative, Nominal (caution in real life application)



DESCRIPTIVE STATISTICS – 1

• Measure of central tendency

- Mean: Measure of the average of all the values in a sample is called Mean
- Median: Measure of the central value of the sample set is called Median. (Is p50 robust to data skewness, is a robust estimator)
- Mode: The value most recurrent in the sample set is known as Mode. (Used frequently for categorical variable, sometime used for discrete numerical)
- Weighted Mean (Two major reasons to use it)

$$\text{Mean} = \bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

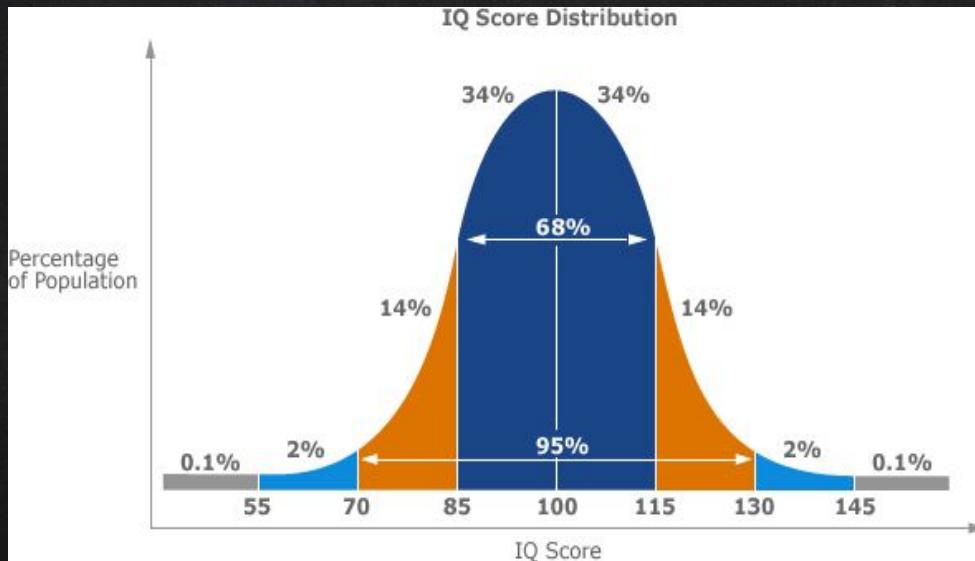
$$\text{Weighted mean} = \bar{x}_w = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$

• Measure of spreadness

- Variance
- Standard deviation
- Quartile: Quartiles tell us about the spread of a data set by breaking the data set into quarters, just like the median breaks it in half (Robust).
- Range: It is the given measure of how spread apart the values in a data set are. The range can be calculated as: $\text{Range} = \text{Max}(x_i) - \text{Min}(x_i)$

$$\begin{aligned}\text{Variance} &= s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \\ \text{Standard deviation} &= s = \sqrt{\text{Variance}}\end{aligned}$$

DESCRIPTIVE STATISTICS – 1 EXAMPLE





DESCRIPTIVE STATISTICS – 1 PRACTICE

Practice: The salaries of staff members is given below:

Case 1: Staff salary = 15k, 18k, 16k, 14k, 15k, 15k, 12k, 17k, 90k, 95k

Mean? Median? Mode? Percentiles?

Mean Median Mode Min P25 P50 P75 Max

Case 2: Staff salary = 15k, 18k, 16k, 14k, 13k, 15k, 17k, 17k, 14k, 14k (bold one's are outliers)

Mean? Median? Mode? Percentile

Mean Median Mode Min P25 P50 P75 Max

$$\text{Mean} = \bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

DESCRIPTIVE STATISTICS – 1 PRACTICE ANSWER

Practice: The salaries of staff members is given below:

Case 1: Staff salary = 15k, 18k, 16k, 14k, 15k, 15k, 12k, 17k, 90k, 95k

Mean? Median? Mode? Percentiles?

Mean	Median	Mode	Min	P25	P50	P75	Max
30.7K	15.5K	15K	12K	15K	15.5K	18K	95K

Case 2: Staff salary = 15k, 18k, 16k, 14k, 13k, 15k, 17k, 17k, 14k, 14k (bold one's are outliers)

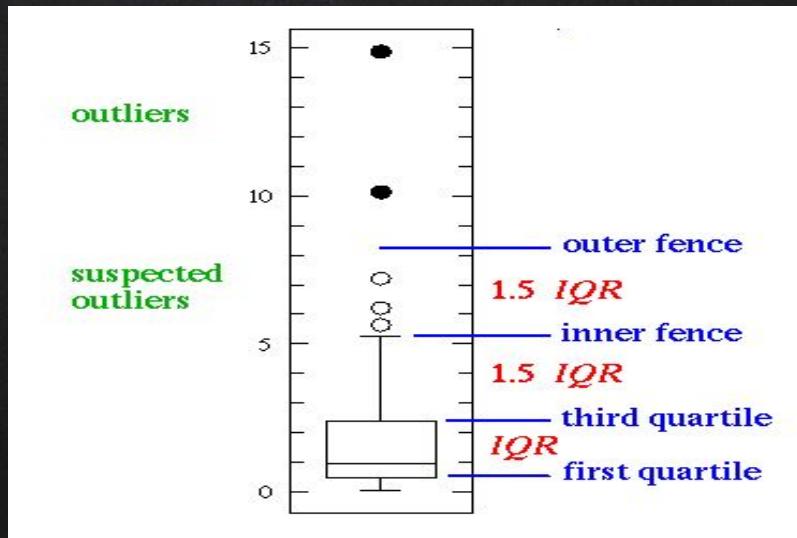
Mean? Median? Mode? Percentile

Mean	Median	Mode	Min	P25	P50	P75	Max
15.3K	15K	14K	13K	14K	15K	17K	18K

$$\text{Mean} = \bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

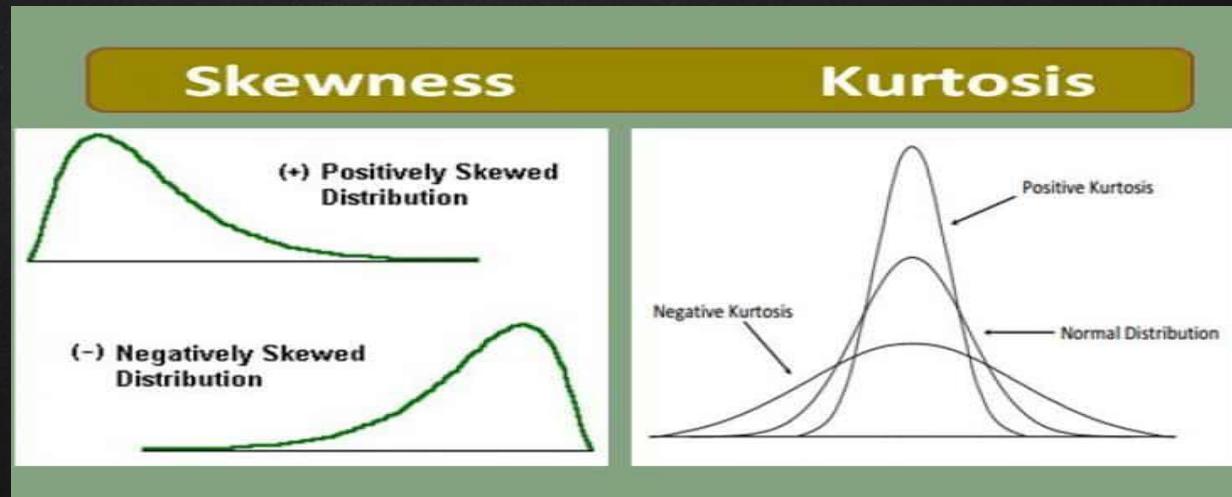
DESCRIPTIVE STATISTICS - 2

- Outlier and IQR(Interquartile Range)
- Five number summary statistics - Boxplot (One of the simplest data distribution)



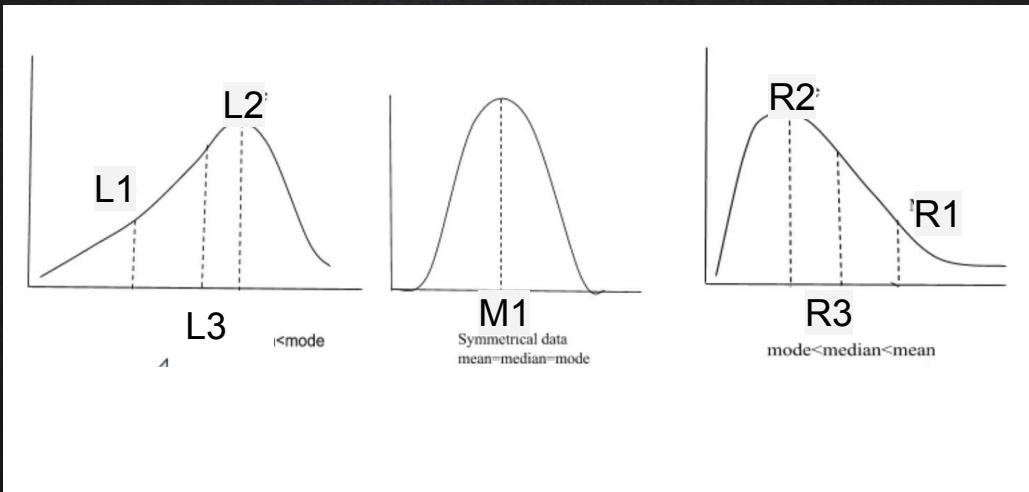
DESCRIPTIVE STATISTICS FOR EXPLORING DATA DISTRIBUTION

- Skewness
- Kurtosis

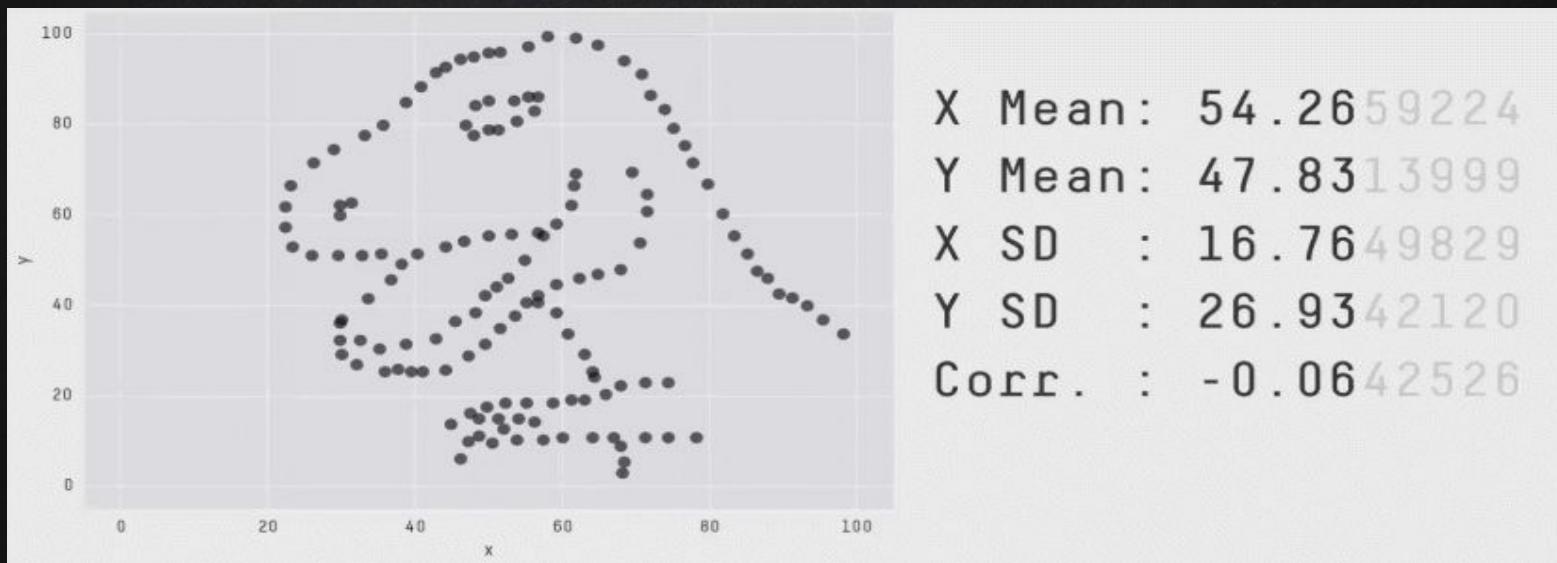


EXPLORING DATA DISTRIBUTION – PRACTICE

Could you label Mean, Median and Mode on the three charts ?

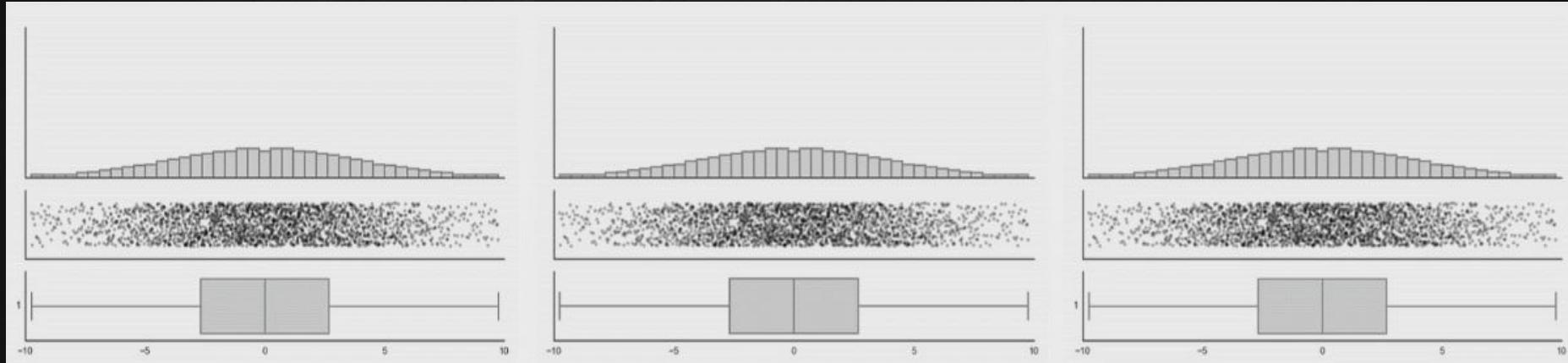


IS MEAN AND VARIANCE SUFFICIENT TO SUMMARIZE/UNDERSTAND YOUR DATA?



<https://www.autodeskresearch.com/publications/samestats>

Is BOXPLOT SUFFICIENT?



<https://www.autodeskresearch.com/publications/samestats>



KEY CHARTS TO EXPLORE DATA DISTRIBUTION – UNIVARIATE

Boxplot

"A plot introduced by Tukey as a quick way to visualize the distribution of data"

Frequency table/Bar Chart

"A tally of the count of categorical data for each category (numerical variable is histogram)"

Histogram

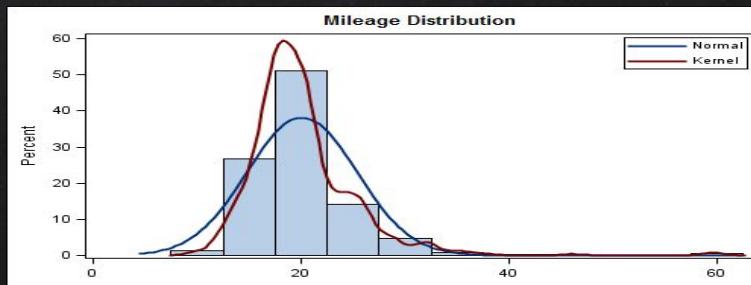
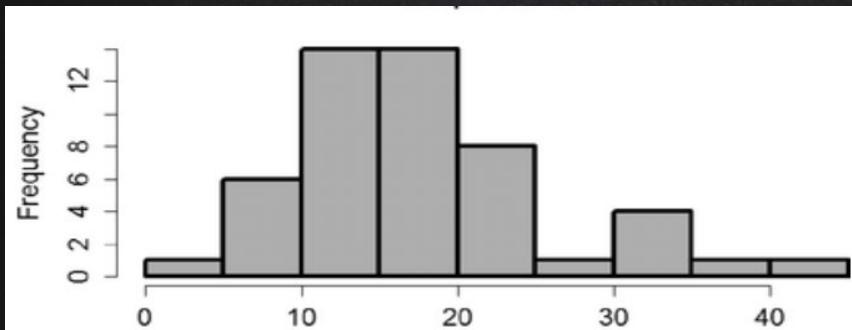
"A plot of the frequency table with the bins on the x-axis and the count (or proportion) on the y-axis. While visually similar, bar charts should not be confused with histograms"

Density plot

"A smoothed version of the histogram, often based on a *kernel density estimate*"



HOW THESE CHARTS LOOKS LIKE





KEY CHARTS TO EXPLORE DATA DISTRIBUTION – BIVARIATE

Scatter plot - correlation (two numerical variables)

"A plot in which the x-axis is the value of one variable, and the y-axis the value of another"

Heatmap - correlation (two numerical variables)

"A plot build on correlation matrix, showed the any two variable correlation in a heatmap color theme"

Correlation matrix

"A table where the variables are shown on both rows and columns, and the cell values are the correlations between the variables"

Contour plot

"A plot showing the density of two numeric variables like a topographical map"

Violin plot

"Similar to a boxplot but showing the density estimate"

Bar plot

"A boxplot for several bars showed one numerical measure across a nominal variable category"

Contingency table

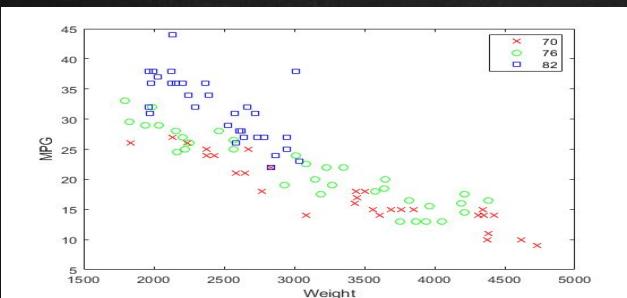
"A tally of counts between two or more categorical variables"

BIVARIATE CHARTS EXAMPLE1

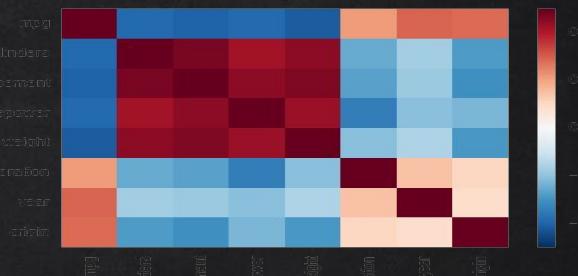
Correlation Matrix

	T	CTL	FTR	VZ	LVLT
T	1.000	0.475	0.328	0.678	0.279
CTL	0.475	1.000	0.420	0.417	0.287
FTR	0.328	0.420	1.000	0.287	0.260
VZ	0.678	0.417	0.287	1.000	0.242
LVLT	0.279	0.287	0.260	0.242	1.000

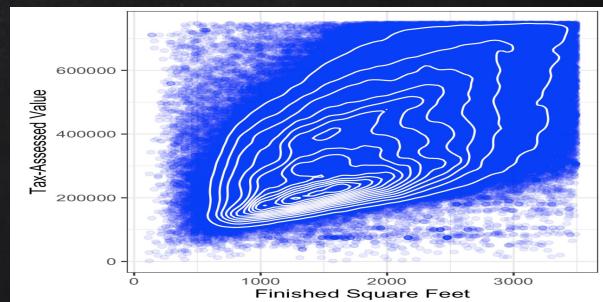
Scatter Plot



Correlation Heatmap



Contour Plot

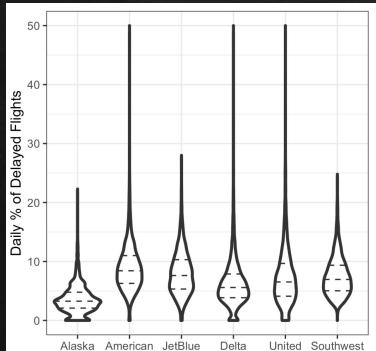


BIVARIATE CHARTS EXAMPLE2

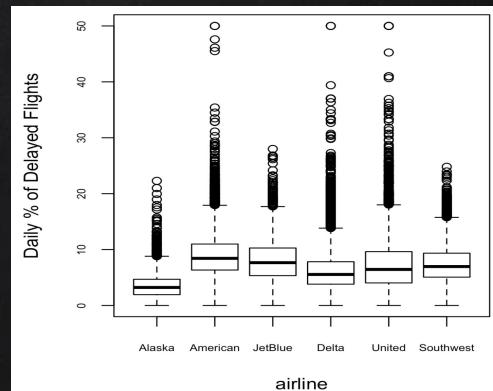
Contingency Table

	Dog	Cat	Total
Male	42	10	52
Female	9	39	48
Total	51	49	100

Violin Plot

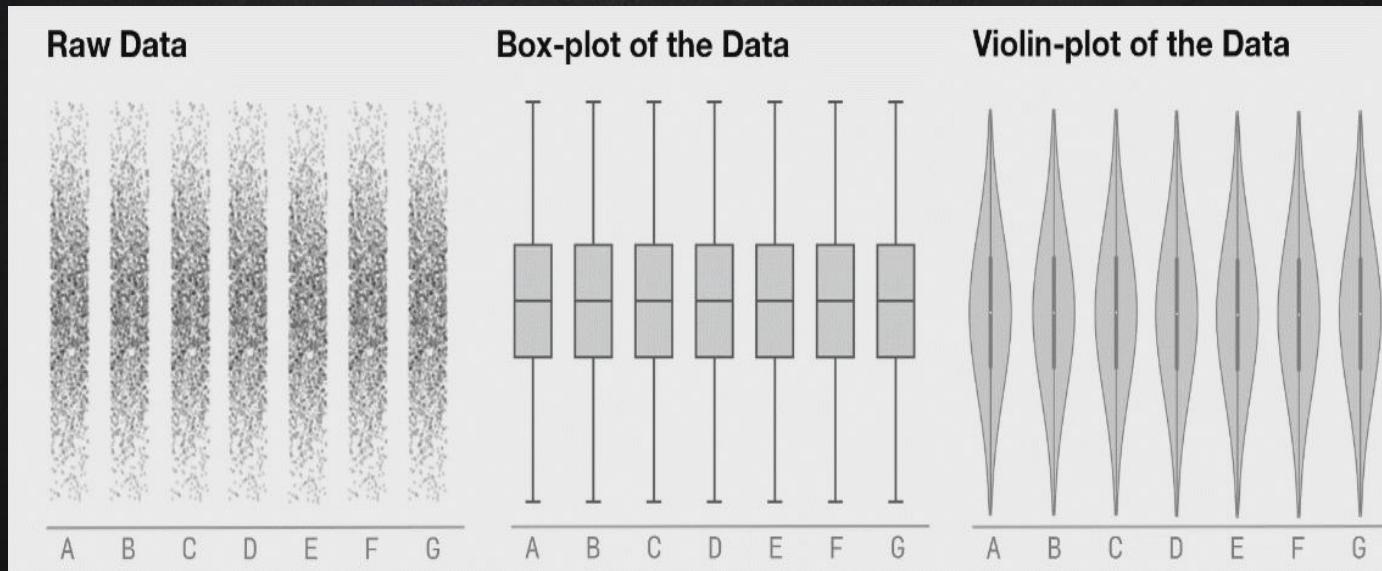


Bar Plot





BIVARIATE BOXPLOT NORMALLY OK ...



<https://www.autodeskresearch.com/publications/samestats>



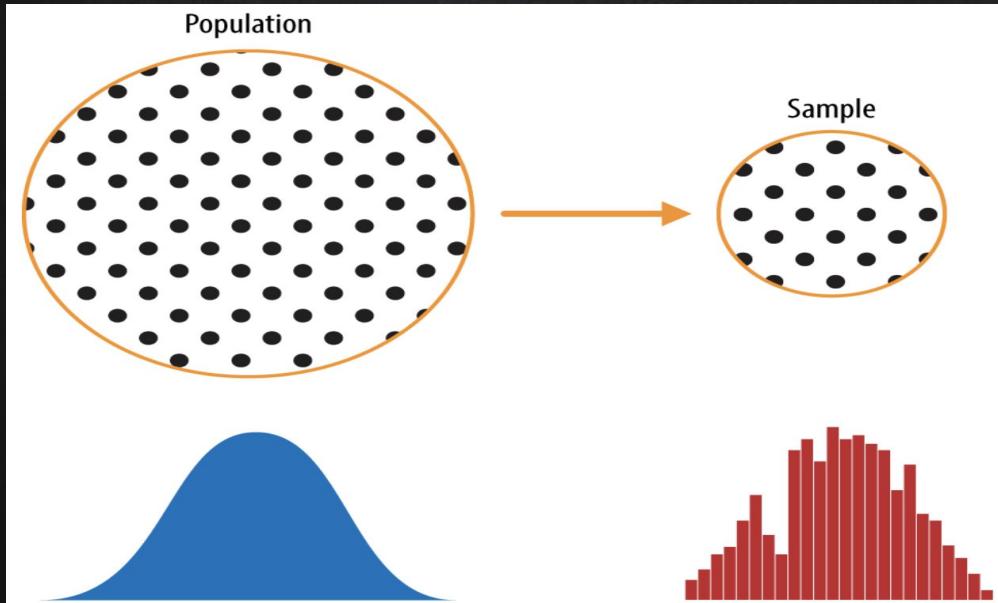
TABLE OF CONTENTS FOR TODAY'S DISCUSSION

1. Basics Statistics to understand data (EDA statistics)

2. Sampling Techniques

3. Probability distribution and Inferential statistical analysis

POPULATION/SAMPLE DISTRIBUTION



“Data quality often matters more than data quantity when making an estimate or a model based on a sample. Data quality in data science involves completeness, consistency of format, cleanliness, and accuracy of individual data points. Which means *representativeness*. ”

SAMPLE BIAS AND SAMPLING

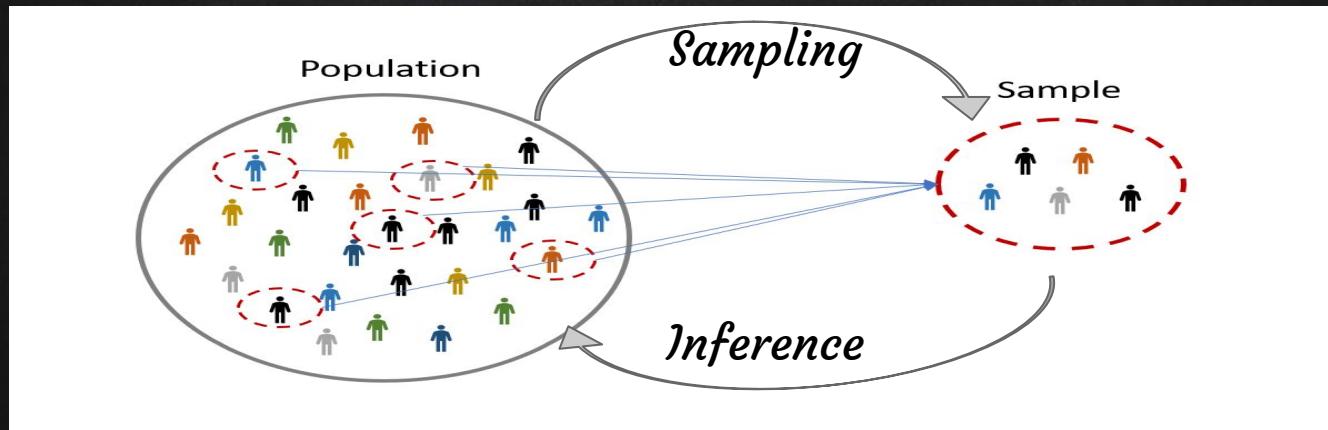
Type of sample bias:

- Selection bias
- Data collection bias
- Prescreen bias
- Non-probability sample bias – unquantifiable

Sampling Designs help overcome bias:

- Random sampling: "Drawing elements into a sample at random"
- Stratified sampling: "Dividing the population into strata and randomly sampling from each strata"
- Simple Random Sample: "The sample was results from random sampling without stratifying"
- Clustered Sample: "occurs when the natural sampling unit is a group or cluster of individual units"
- Systematic Sample: "the selection of every kth element from a sampling frame or from a sequential stream of potential respondents"

SAMPLE STATISTICS ARE RANDOM VARIABLES WITH DISTRIBUTION



Population parameters:
population mean, population
variance

Sample statistics: sample mean,
sample variance

SAMPLING DISTRIBUTION

* Sampling Distribution

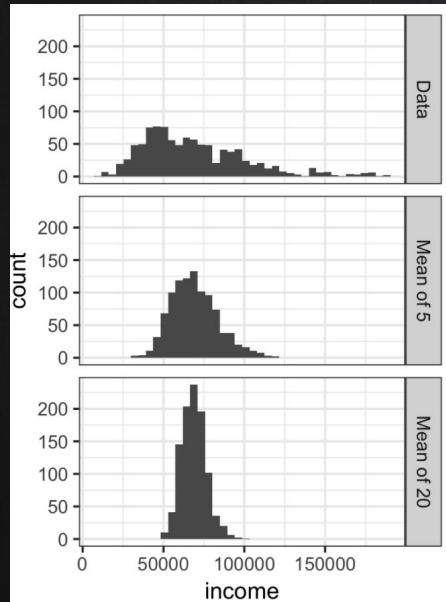
Parameter
Random sample of size n → statistic
 $\mu = \frac{1+2+3}{3} = 2$

# s pick	\bar{x}
1, 1	1
1, 2	1.5
1, 3	2
2, 1	1.5
2, 2	2
2, 3	2.5
3, 1	2
3, 2	2.5
3, 3	3

Khan Academy

[Sampling distribution video from Khan Academy](https://www.youtube.com/watch?v=z0Ry_3_qhDw)
https://www.youtube.com/watch?v=z0Ry_3_qhDw

DATA DISTRIBUTION AND SAMPLING DISTRIBUTION



Sample statistic

"A metric calculated for a sample of data drawn from a larger population"

Data distribution

"The frequency distribution of individual *values* in a data set"

Sampling distribution

"The frequency distribution of a *sample statistic* over many samples or resamples"

Central limit theorem

"The tendency of the sampling distribution to take on a normal shape as sample size rises"



STANDARD ERROR OF THE MEAN

The thumbnail features a purple bell curve icon on the left. In the center, the title 'Standard error of the mean' is written in white, bold, sans-serif font. To the left of the title is a green stylized letter 'σ'. Below the title, the 'Khan Academy' logo is visible, consisting of a green hexagon with a white silhouette of a person inside.

Mean and standard error of sample mean video from Khan Academy

<https://www.youtube.com/watch?v=J1twbrHeI3o>

RESAMPLING METHODS

Bootstrap: Random sample with replacement

Data: A, B, C, D

Samples: (A/B), (B/C), (A/B), (C/D).....

Permutation: Random sample without replacement

Data: A, B, C, D

Samples: (A/B), (A/C), (A/D), (B/C), (B/D), (C/D)

Bootstrap is the most popular resampling method (an non-parametric approach) for parameter estimate as well as some non-parametric ML algorithms

Warning: The bootstrap does not compensate for a small sample size; it does not create new data, nor does it fill in holes in an existing data set. It merely informs us about how lots of additional samples would behave when drawn from a population like our original sample



TABLE OF CONTENTS FOR TODAY'S DISCUSSION

1. Basics Statistics to understand data (EDA statistics)

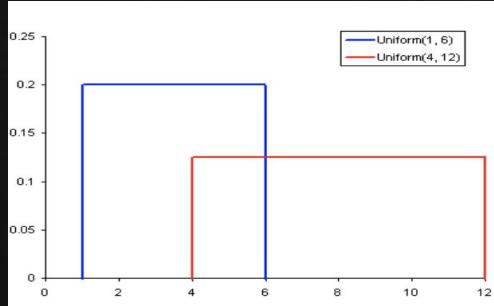
2. Sampling Techniques

3. Probability distribution and Inferential statistical analysis

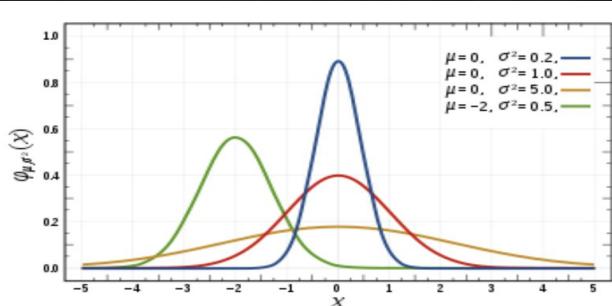
PROBABILITY DENSITY FUNCTION

Probability is the measure of how likely an event will occur. More precisely, probability is the ratio of desired outcomes to total outcomes: $(\text{desired outcomes}) / (\text{total outcomes})$

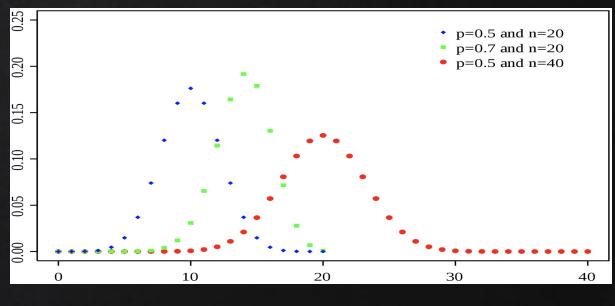
Uniform



Normal



Binomial



EXPERIMENT AND INFERRENTIAL STATISTICAL ANALYSIS

A/B Testing: "An A/B test is an experiment with two groups to establish which of two treatments, products, procedures, or the like is superior. Often one of the two treatments is the standard existing treatment, or no treatment. If a standard (or no) treatment is used, it is called the control. A typical hypothesis is that a new treatment is better than the control"

Examples are common in web design and marketing, drug effectiveness

- Testing two soil treatments to determine which produces better seed germination
- Testing two therapies to determine which suppresses cancer more effectively
- Testing two prices to determine which yields more net profit
- Testing two email headlines to determine which produces more clicks
- Testing two web ads to determine which generates more conversions

STATISTICAL TESTING

The Null Hypothesis: treatment A has the same response as B

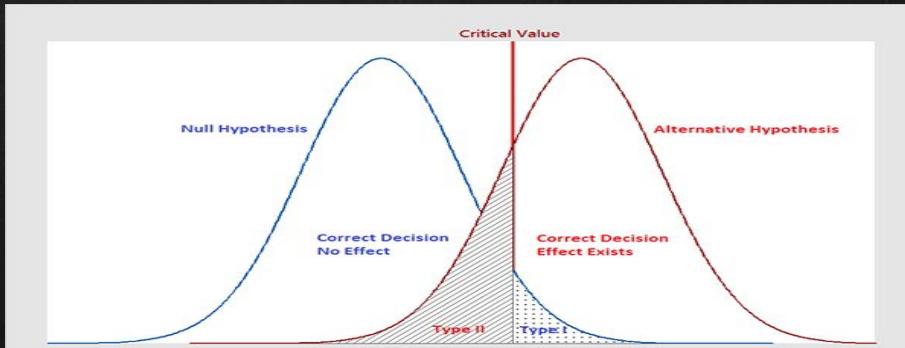
Alternative Hypothesis: treatment A has better response than B

P-value: Under the null hypothesis (Or if the stated hypothesis is true), the p-value is the probability of obtaining a result(statistic) as unusual or extreme as the observed results.

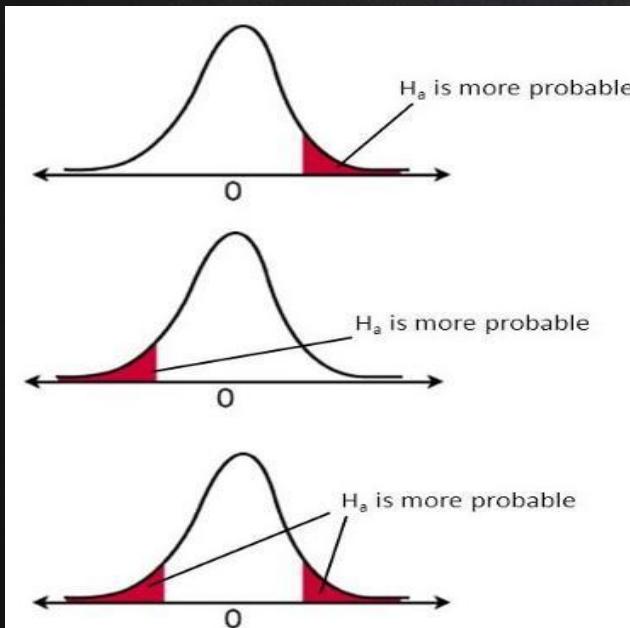
Alpha - significance level: The probability threshold of “unusualness” that chance results must surpass for actual outcomes to be deemed statistically significant. Or Probability of making a Type I Error: Reject null hypothesis when null hypothesis is true

Type 1 error: error occurs when the sample results, lead to the rejection of the null hypothesis when it is in fact true. Type-I errors are equivalent to false positives.

Type 2 error: error occurs when based on the sample results, the null hypothesis is not rejected when it is in fact false. Type-II errors are equivalent to false negatives.



ONE-TAILED AND TWO-TAILED TESTS



Right-tail test

$$H_a: \mu > \text{value}$$

Left-tail test

$$H_a: \mu < \text{value}$$

Two-tail test

$$H_a: \mu \neq \text{value}$$

One Tailed Test

Two Tailed Test



TEST STATISTIC:

"The **test statistic** measures how close the sample has come to the null hypothesis. Its observed value changes randomly from one random sample to a different sample. A test statistic contains information about the data that is relevant for deciding whether to reject the null hypothesis or not"

One sample test and test statistic

Hypothesis Test	Test Statistics	Expression
Z-Test	Z-statistics	$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$
T-test	t-statistics	$\frac{\bar{X} - \mu}{S/\sqrt{n}}$

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$



THE DIFFERENCE BETWEEN Z-STATISTICS AND T-STATISTICS

Z $\approx \frac{\bar{x} - \mu_{\bar{x}}}{s / \sqrt{n}}$ "OK" if $n > 30$

Z-statistics vs. T-statistics

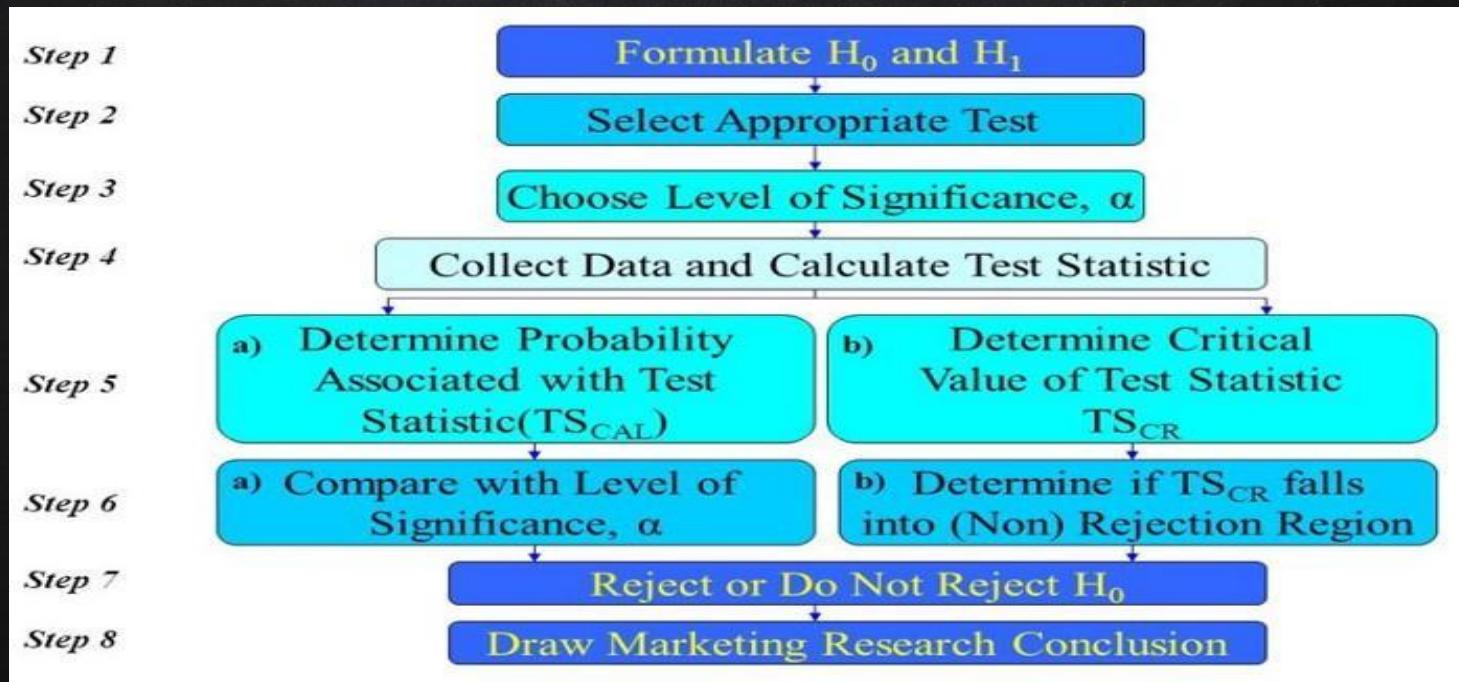
t-distribution $\frac{s}{\sqrt{n}}$ small

Khan Academy

T and Z statistics video from Khan Academy
Difference of T and Z statistics



THE NORMAL PROCEDURE FOR AB TESTING





EXERCISES

A major department store is considering the introduction of an Internet shopping service. The new service will be introduced if more than 40 percent of the Internet users shop via the Internet.

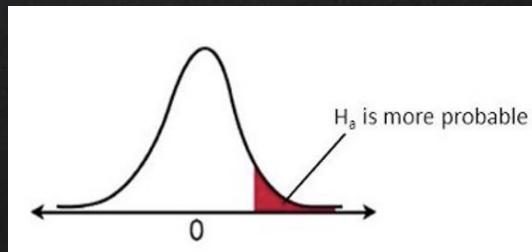
The experiment result: 30 users were surveyed and 17 indicated that they used the Internet for shopping.

Please following the procedure to conduct a hypothesis testing to suggest - should the department store introduce the internet shopping service?

SOLUTION – STEP1 FORMULATING THE HYPOTHESIS

$$H_0: \pi \leq 0.40$$

$$H_1: \pi > 0.40$$





SOLUTION – STEP2 SELECT A SUITABLE TEST

This is a proportion problem, binomial distribution and follow the central limit theorem, when n is large enough, the sample mean following normal distribution, then Z test and z-statistic will be the one to use

$$Z = (p - \pi) / \sigma \text{ where } \sigma = \sqrt{\pi(1-\pi)/n}$$



SOLUTION – STEP3 CHOOSE LEVEL OF SIGNIFICANCE, α

Level of Significance refers to **Type-I** error. In this problem, a Type-I error would occur if we concluded, based on the sample data, that the proportion of customers preferring the new service plan was greater than 0.40, when in fact it was less than or equal to 0.40.

On the other hand, **Type-II** error should occur when the sample data help concluded that the proportion of customer would prefer the new service is less than or equal to 0.4, while the true proportion is larger than 0.4.

The choice should balance the two types of error, normally 0.05 and 0.01 are commonly used significance levels, other choice are rare

Choice $\alpha = 0.05$



SOLUTION – STEP4: COLLECT DATA AND CALCULATE TEST STATISTICS

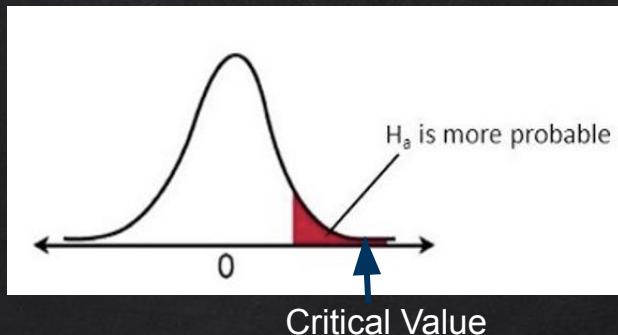
N=30, Yes: 17

Therefore $p = 17/30 = 0.567$

For $\sigma = \sqrt{0.4 * (1-0.4)/30} = 0.089$

$$\begin{aligned}\text{Finally: } z &= (p - \pi) / \sigma \\ &= (0.567 - 0.4) / 0.089 \\ &= 1.88\end{aligned}$$

SOLUTION – STEP5 DETERMINE THE TEST STATISTIC PROBABILITY(P-VALUE)



Approach 1:

$$\Pr(z \leq 1.88) = 0.9699$$

$$P\text{-value} = 1 - \Pr(z \leq 1.88) = 0.03$$

Approach 2:

Critical Value(CV) (z when $\Pr(z \leq CV) = 0.95$)

$$CV = 1.65$$



SOLUTION – STEP 6&7 COMPARE WITH SIGNIFICANT LEVEL AND DRAW CONCLUSION

$$P\text{-value} = 0.03 < 0.05 = \alpha \quad \text{OR} \quad CV = 1.65 < 1.88$$

Reject the Null hypothesis, accept the alternative hypothesis,
the customer proportion of prefer the new service is
greater than 0.4



SOLUTION – STEP 8 CONCLUSION AND MORE

THIS TEST CONCLUDED THAT THERE IS EVIDENCE THAT THE PROPORTION OF INTERNET USERS WHO SHOP VIA THE INTERNET IS SIGNIFICANTLY GREATER THAN 0.40. HENCE, THE RECOMMENDATION TO THE DEPARTMENT STORE WOULD BE TO INTRODUCE THE NEW INTERNET SHOPPING SERVICE.



ASSUMPTION OF ONE-SAMPLE PARAMETRIC STATISTICAL TESTING

- One-Sample Z-Test Assumptions:
 1. The data are continuous (not discrete).
 2. The data follow the normal probability distribution.
 3. The sample is a simple random sample from its population. Each individual in the population has an equal probability of being selected in the sample.
 4. The population standard deviation is known.
 5. Binary data proportion test validated by Central Limit Theorem for normality, with $n*\pi > 5$ and $n*(1-\pi) > 5$
- One-Sample T-Test Assumptions The assumptions of the one-sample t-test are:
 1. The data are continuous (not discrete).
 2. The data follow the normal probability distribution.
 3. The sample is a simple random sample from its population. Each individual in the population has an equal probability of being selected in the sample.



TYPES OF STATISTICAL TESTING I

- One sample test

T-Test and Z-test

- Two sample test with equal variance

T-Test and Z-test

- Two sample test with unequal variance

T-Test and Z-test

Detail information

Hypothesis Test	Test Statistics	Expression
One sample Z-Test	Z-statistics	$\frac{\bar{X} - \mu}{\sigma / \sqrt{n}}$
One sample T-test	T-statistics	$\frac{\bar{X} - \mu}{S / \sqrt{n}}$
Two sample Z-Test equal variance	Z-statistics	$z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$
Two sample Z-Test unequal variance	Z-statistics	$z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$
Two sample T-Test equal variance	T-statistics	$T = \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}}$
Two sample T-Test unequal variance	T-statistics	$t = \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{\sqrt{\frac{S_X^2}{n} + \frac{S_Y^2}{m}}}$

delta1=delta2



TYPES OF STATISTICAL TESTING II

Parametric: "assume underlying statistical distributions in the data. Therefore, several conditions of validity must be met so that the result of a parametric test is reliable. For example, Student's t-test for two independent samples is reliable only if each sample follows a normal distribution and if sample variances are homogeneous."

Nonparametric: "do not rely on any distribution. They can thus be applied even if parametric conditions of validity are not met."

Pros and cons:

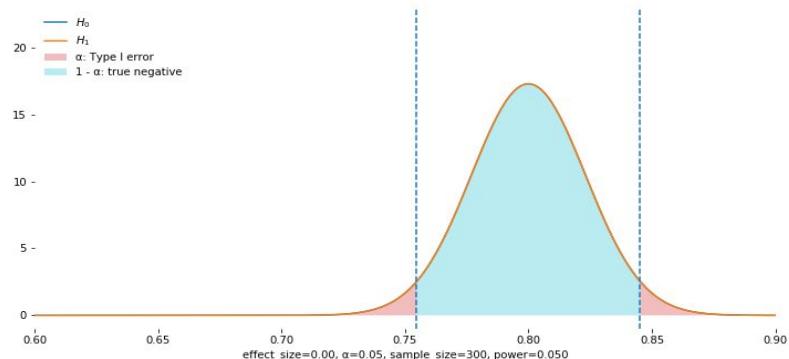
- "Nonparametric tests are more **robust** than parametric tests. In other words, they are valid in a broader range of situations (fewer conditions of validity)."
- "The advantage of using a parametric test instead of a nonparametric equivalent is that the former will have more statistical **power** than the latter. In other words, a parametric test is more able to lead to a rejection of H₀. Most of the time, the p-value associated to a parametric test will be lower than the p-value associated to a nonparametric equivalent that is run on the same data."

POWER AND SAMPLE SIZE

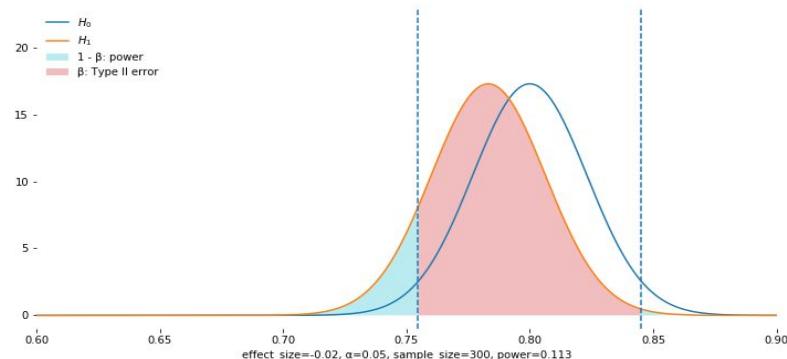
Before we begin the data collection, we should have put several things in consideration beforehand to make the most out of a experimental design to fulfill the purpose of our study to prevent:

1. Test was not sensitive because the sample size was not large enough
2. Is the test provide enough power to test the difference if the experiment is truly different

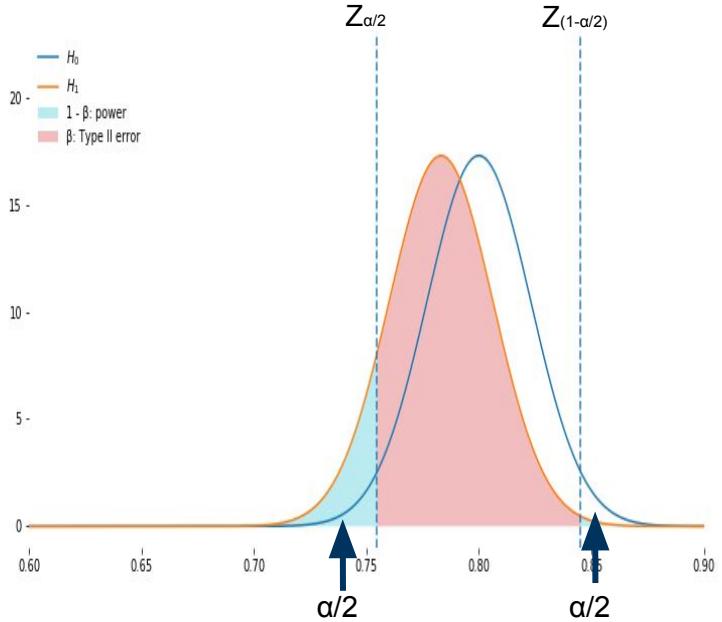
When Null Hypothesis is True



When Null Hypothesis is False



How To CALCULATE POWER



$$H_0: \mu = \mu_0 \quad H_a: \mu \neq \mu_0$$

$$\text{Test statistic: } Z = \frac{(\bar{y} - \mu)}{\sigma/\sqrt{n}} = \frac{\sqrt{n}(\bar{y} - \mu)}{\sigma}$$

significance level α :

$$\alpha = \Pr\{\text{Reject } H_0 / H_0 \text{ True}\} = \Pr\{Z < Z_{\frac{\alpha}{2}} \text{ or } Z > Z_{(1-\frac{\alpha}{2})} / \mu = \mu_0\}$$

Power:

$$1 - \beta = \Pr\{\text{Reject } H_0 / H_0 \text{ False}\} = \Pr\{Z < Z_{\frac{\alpha}{2}} \text{ or } Z > Z_{(1-\frac{\alpha}{2})} / \mu \neq \mu_0\}$$

$$= \Pr\{Z < Z_{\frac{\alpha}{2}} / \mu \neq \mu_0\} + \Pr\{Z > Z_{(1-\frac{\alpha}{2})} / \mu \neq \mu_0\}$$

Assume: H_0 with mean μ_0

H_a with mean μ_1 .

$$\begin{aligned} \text{Power} &= 1 - \beta = \Pr\left\{\frac{\sqrt{n}(\bar{y} - \mu_0)}{\sigma} < Z_{\frac{\alpha}{2}} / \mu = \mu_1\right\} + \Pr\left\{\frac{\sqrt{n}(\bar{y} - \mu_0)}{\sigma} > Z_{(1-\frac{\alpha}{2})} / \mu = \mu_1\right\} \\ &= \Pr\left\{\frac{\sqrt{n}(\bar{y} - \mu_0 + \mu_1 - \mu_0)}{\sigma} < Z_{\frac{\alpha}{2}} / \mu = \mu_1\right\} + \Pr\left\{\frac{\sqrt{n}(\bar{y} - \mu_0 + \mu_1 - \mu_0)}{\sigma} > Z_{(1-\frac{\alpha}{2})} / \mu = \mu_1\right\} \\ &= \Pr\left\{\frac{\sqrt{n}(\bar{y} - \mu_1)}{\sigma} < Z_{\frac{\alpha}{2}} - \frac{\sqrt{n}(\mu_0 - \mu_1)}{\sigma} / \mu = \mu_1\right\} + \Pr\left\{\frac{\sqrt{n}(\bar{y} - \mu_1)}{\sigma} > Z_{(1-\frac{\alpha}{2})} - \frac{\sqrt{n}(\mu_0 - \mu_1)}{\sigma} / \mu = \mu_1\right\} \end{aligned}$$



WHAT'S THE CONSTRAINTS OF SIZE AND POWER CALCULATION

TO CALCULATE POWER OR REQUIRED SAMPLE SIZE, THERE ARE FOUR MOVING PARTS:

SAMPLE SIZE

EFFECT SIZE YOU WANT TO DETECT

SIGNIFICANCE LEVEL (ALPHA) AT WHICH THE TEST WILL BE CONDUCTED

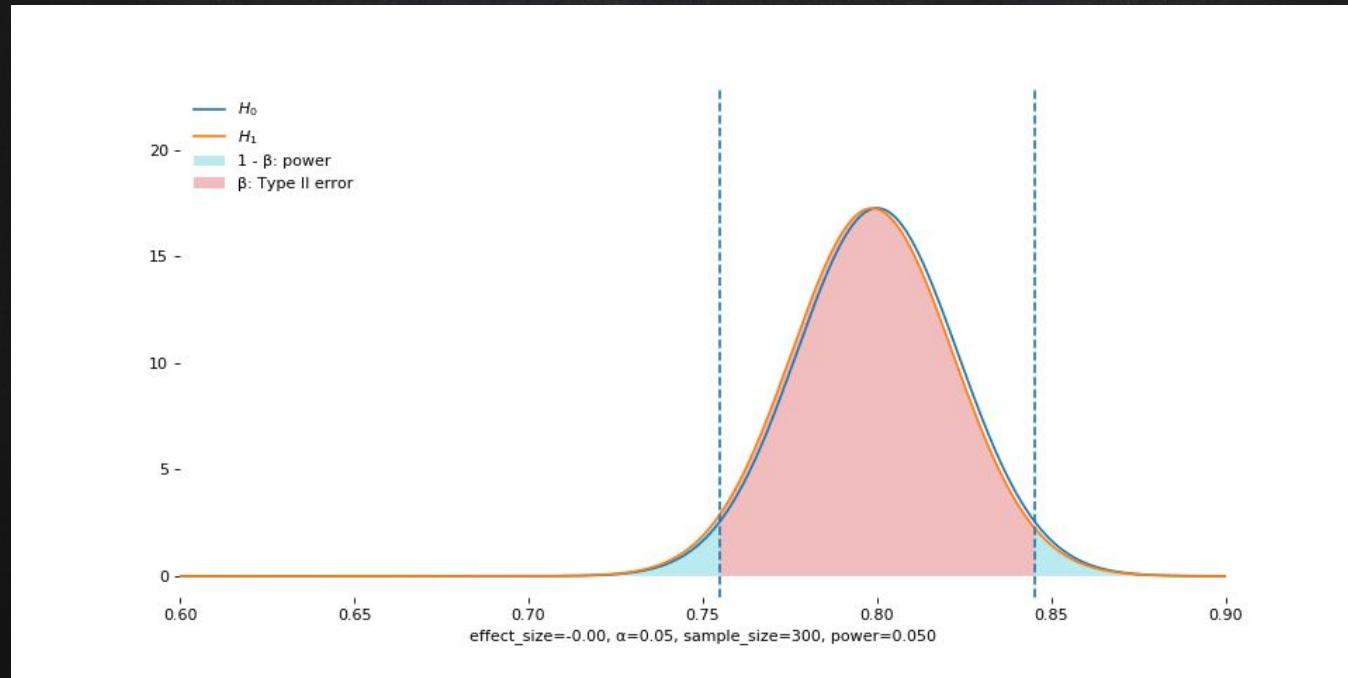
POWER

SPECIFY ANY THREE OF THEM, AND THE FOURTH CAN BE CALCULATED. MOST COMMONLY, YOU WOULD WANT TO CALCULATE SAMPLE SIZE, SO YOU MUST SPECIFY THE OTHER THREE.

$$z = (p - p_0) / \sqrt{(p_0(1-p_0)/n)}$$

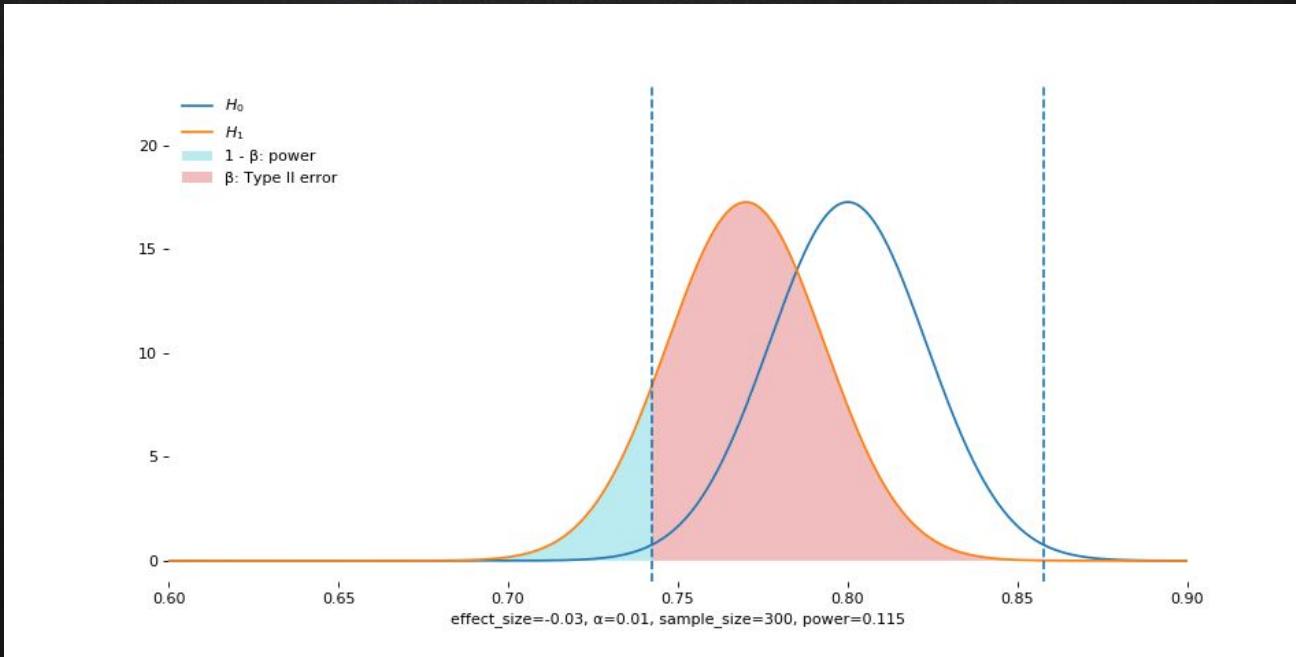


How EFFECTIVE SIZE IMPACT POWER





HOW ALPHA IMPACT POWER



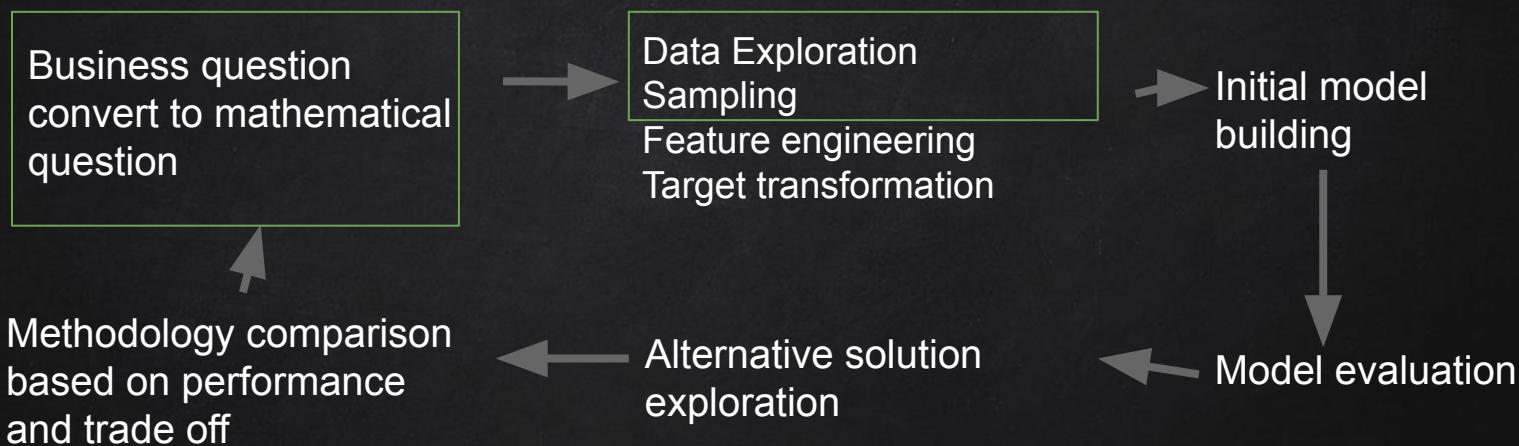


How SAMPLE SIZE IMPACT POWER





GENERAL PROCESS FOR DATA SCIENCE PROJECT



REFERENCES



<https://thefactfactor.com/facts/management/statistics/variable/1432/>

<https://towardsdatascience.com/8-fundamental-statistical-concepts-for-data-science-9b4e8a0c6f1c>

<https://learning.oreilly.com/library/view/Practical+Statistics+for+Data+Scientists,+2nd+Edition/9781492072935/ch03.html#Significance>

https://en.wikipedia.org/wiki/Sampling_bias

<https://towardsdatascience.com/everything-you-need-to-know-about-hypothesis-testing-part-i-4de9abebbc8a>

<https://towardsdatascience.com/the-power-of-a-b-testing-3387c04a14e3>

<https://www.edureka.co/blog/statistics-and-probability/#Inferential%20Statistics>

<https://towardsdatascience.com/6-different-ways-to-compensate-for-missing-values-data-imputation-with-examples-6022d9ca0779>

<https://towardsdatascience.com/the-5-basic-statistics-concepts-data-scientists-need-to-know-2c96740377ae>